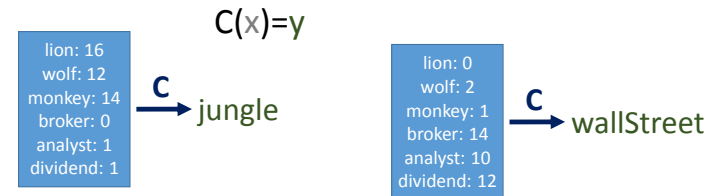


Learning Theory

CISC 5800
Professor Daniel Leeds

The classifier

Function C that provides
correct label (Y) based on features (X)



Goal: identifier classifier that maximizes
correct labels for most inputs

2

Sample complexity

How many training examples needed to learn concept?

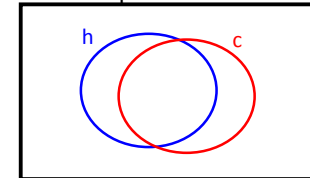
- X – set of data points
- $P(X)$ – Probability of drawing data point x
- H – space of hypotheses $H = \{h : X \rightarrow \text{classes}\}$
- C – correct assignment $C = \{c : c(x) = y \forall x \in X\}$

3

Probability of error

$$H = \{h : X \rightarrow \{0,1\}\}$$

X – data points



True error of h : probability randomly selected
data point from $P(X)$ misclassified

$$\text{error}_{\text{true}}(h) = \Pr_{x \sim P(X)} [h(x) \neq c(x)]$$

- Hard to compute, but can prove properties of $\text{error}_{\text{true}}$

4

Example: Learner picks one of fixed number of classifiers $h \in H$

Correct classifier c is some assignment of each x to a label

How many training points m needed for $\text{error}_{\text{true}}(h) < \varepsilon$?

$$\text{Prob}[\text{error}_{\text{true}}(h) \leq \varepsilon] > 1 - \delta$$

“Probability learned classifier h has worse than ε error is $< \delta$ ”

“Probably Approximately Correct Learning” – PAC Learning

5

Binary example: sample complexity

Note for $\varepsilon \in [0,1]$, $(1 - \varepsilon) \leq e^{-\varepsilon}$

What is the chance learned h is bad but classifies training data correctly?

If $\text{error}_{\text{true}}(h) > \varepsilon$:

- $\text{Prob}[h \text{ correctly labels } x^1] < (1 - \varepsilon) \leq e^{-\varepsilon}$
- $\text{Prob}[h \text{ correctly labels } x^1 \text{ and } x^2 \dots \text{ and } x^m] < (1 - \varepsilon)^m \leq e^{-m\varepsilon}$

If classifier picks one h^* randomly from H

- $\text{Prob}[h^* \text{ is bad}] = \text{Prob}[h_1 \text{ bad}] + \dots + \text{Prob}[h_n \text{ bad}]$
 $= \text{Prob}[\text{error}_{\text{true}}(h^*) > \varepsilon] < |H| e^{-m\varepsilon}$ Valiant, 1984

6

Binary example: sample complexity

Number of data points to reduce chance of false classification, enforce

$$\text{Prob}[\text{error}_{\text{true}}(h) \leq \varepsilon] > 1 - \delta$$

$$1 - \text{Prob}[\text{error}_{\text{true}}(h) \leq \varepsilon] = \text{Prob}[\text{error}_{\text{true}}(h) > \varepsilon] < \delta$$

$$\text{Prob}[\text{error}_{\text{true}}(h^*) > \varepsilon] < |H| e^{-m\varepsilon}; \text{ stricter bound } |H| e^{-m\varepsilon} < \delta$$

Valiant, 1984

8

Binary example: sample complexity

Number of data points to reduce chance of false classification, enforce

$$\text{Prob}[\text{error}_{\text{true}}(h) \leq \varepsilon] > 1 - \delta$$

$$\text{Prob}[\text{error}_{\text{true}}(h^*) > \varepsilon] < |H| e^{-m\varepsilon} < \delta$$

$$m > \frac{1}{\varepsilon} \ln \frac{|H|}{\delta}$$

Valiant, 1984

9

VC Dimensions

If H not finite, PAC result seems to require ∞ data points

- Overly conservative

“Dichotomy” – division of set of points S into two subsets

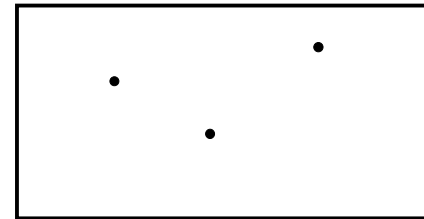
- “Shattering” – set of points is **shattered** by H iff there exists $h \in H$ associated with every possible dichotomy

Vapnik-Chervonenkis dimension **VC(H)** is size of largest finite subset of S that can be shattered by H

10

Shattering example3

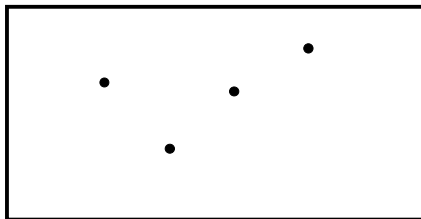
- $H = \{\text{rectangles: inside is 1, outside is 0}\}$ **VC(3)**
- $S = \{3 \text{ specified dots}\}$



12

Shattering example

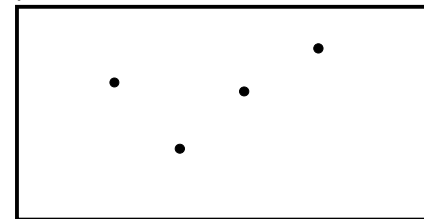
- $H = \{\text{rectangles, inside is 1 outside is 0}\}$ **VC(3)**
- $S = \{4 \text{ specified dots}\}$



14

Shattering example

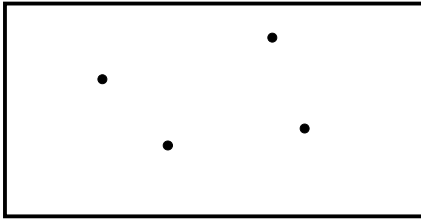
- $H_2 = \{\text{rectangles, inside is 1 outside is 0 inside is 0 outside is 1}\}$ **VC(4)**
- $S = \{4 \text{ specified dots}\}$



15

Shattering example

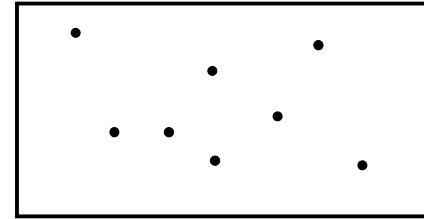
- $H = \{\text{rectangles, inside is 1 outside is 0}\}$ **VC(4)**
- $S = \{4 \text{ specified dots}\}$



18

Shattering example 4

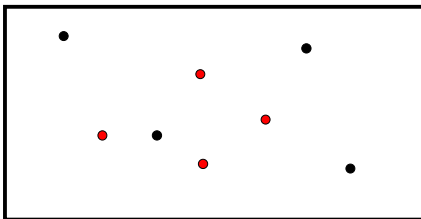
- $H = \{\text{rectangle, inside is 1 outside is 0}\}$
- $S = \{8 \text{ specified dots}\}$



21

Shattering example 4

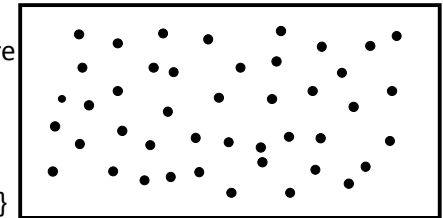
- $H = \{\text{rectangle, inside is 1 outside is 0}\}$ **H(4)**
- $S = \{8 \text{ specified dots}\}$



22

Shattering infinite points

- $H = \{\text{Linear separators}\}$
- $S = \{\text{Any point in 2D feature space}\}$
- $S = \{\text{Any point in } nD \text{ space}\}$



23

PAC result with infinite H

VC(H) is size of largest finite subset of X that can be shattered by H

- $d = VC(H)$
- $m \geq O\left(\frac{1}{\epsilon} \left[d \log \frac{1}{\epsilon} + \log \frac{1}{\delta} \right]\right) \sim \frac{1}{\epsilon} \left[d \log \frac{1}{\epsilon} + \log \frac{1}{\delta} \right]$

Recall: $m > \frac{1}{\epsilon} \ln \frac{|H|}{\delta}$ for finite size H

24

Intuition for PAC result with infinite H

- $d = VC(H)$
- $m \geq O\left(\frac{1}{\epsilon} \left[d \log \frac{1}{\epsilon} + \log \frac{1}{\delta} \right]\right) \sim \frac{1}{\epsilon} \left[d \log \frac{1}{\epsilon} + \log \frac{1}{\delta} \right]$
- Finite H: $m > \frac{1}{\epsilon} \ln \frac{|H|}{\delta}$

$$d \log \frac{k}{\epsilon} \rightarrow \log \frac{k^d}{\epsilon}$$

Can pick h to shatter at most d points in one of two classes
 2^d meaningfully different classifiers h: $|H| \sim 2^d$

25